

# A quasi-alternating Markov-modulated linear regression: model implementation using data about coaches' delay time

Nadezda Spiridovska

**Abstract**— This research presents a case-study of quasi-alternating Markov-modulated linear regression application for analysis of delays of coaches (regional buses) on the route Ventspils-Riga in Latvia. Markov-modulated linear regression suggests that the parameters of regression model vary randomly in accordance with external environment which is described as a continuous-time homogeneous irreducible Markov chain with known parameters. Markov-modulated linear regression model differs from other switching models by a new analytical approach. For each state of the environment the regression model parameters are estimated. External environment has only two states in this research that is why model is called quasi-alternating. Data on weather conditions provided by the Latvian Environment, Geology and Meteorology Centre and is free downloaded from its database. Data on weather conditions in the Ventspils city are used for the environment description: two alternate states are assumed: “no precipitation” and “precipitation”. The model of the external environment is tested for the markovian properties using inferential statistics. Actual data on coaches' trip times is provided by the Riga International Coach Terminal. Data are analysed by means of descriptive statistics. Different experiments are carried out and the application of Markov-modulated linear regression model on given sample showed adequate results indicating the validity of the model

**Keywords**— External environment, delay time analysis, Markov-modulated linear regression, trip time.

## I. INTRODUCTION

GENERALLY application of probabilistic-statistical models presupposes invariability of parameters throughout the process of model consideration. In this case it refers to the regression model parameters, i.e., the regression coefficients. However, in practice these parameters usually vary randomly attesting to “random environment” in which investigated object is constantly changing.

Markov-modulated linear regression suggests that the parameters of regression model vary randomly in accordance with the external environment. The latter is described as a continuous-time homogeneous irreducible Markov chain with known parameters. In the classical definition, the term “alternating” refers to the theory of renewal processes, where alternating renewal process or zero-one renewal process has certain theoretical definition. In our case term “alternating” is

used to describe the external environment, which has only two states that change each other or alternate. Thus, the term “alternating” is used in its direct meaning and to avoid a confusion with the classical mathematical definition of the term “alternating” the concept “quasi” is used.

Previous investigations [1, 2] were executed on artificial data (simulation analysis) and showed that in case of small sample estimated parameters, they considerably deviated from true ones (what was explained by insufficient sample size and big randomness of the external environment) and in case of big sample the estimated parameters were very close to true ones, but still convergence to true values was very slow.

Due to cooperation with the Riga International Coach Terminal (RICT), it became possible to put the proposed quasi-alternating Markov-modulated linear regression model into practice using real data. The RICT is a leader in the area of passenger bus transportation services in Latvia (151803 routes per 2015 year). RICT serves approximately more than 1.860 million of passengers per year [4]. Punctuality and accurate adherence to a timetable are ones of the most significant factors affecting the services quality level. RICT management annually conducts punctuality analysis as part of quality management system. In this research coaches delay time on the route Ventspils-Riga is analysed.

Section 2 provides detailed description of a quasi-alternating Markov-modulated linear regression. Section 3 presents data chosen for the external environment description. Section 4 contains data description of coaches' delay time. And section 5 presents modelling results: implementation of Markov-modulated linear regression using data about coaches' delay time.

## II. A QUASI-ALTERNATING MARKOV-MODULATED LINEAR REGRESSION: THE OUTLINE

The idea of combining in particular way linear regression models and Markov-chain based models was put forward by professor Alexander Andronov and was first described in [1,2].

Markov-Modulated linear regression model assumes that the external environment is described by an irreducible Markov chain with continuous time, parameters of which are known. Let us consider the main idea of the model.

We suppose that the classical regression model of the form  $Y_i = x_i \beta + Z_i$ ,  $i=1, \dots, n$ , (where  $Y_i$  is scale response,  $x_i$  is a vector,  $\beta$  is a vector,  $Z_i$  is a scale disturbance,  $n$  is a number of observations) corresponds to one unit of a continuous time,  $Z_i$

The author is with the Research Department, Mathematical Methods and Modelling Department, Transport and Telecommunication Institute Lomonosova 1, Riga, LV-1019, LATVIA, Spiridovska.N@tsi.lv

is Brown motion and responses  $Y_i(t)$  are time-additive. Then for  $t > 0$ :

$$Y_i = x_i \beta t + Z_i \sqrt{t}, i = 1, \dots, n, \quad (1)$$

where  $Y_i(0) = 0$ . Further value  $Y_i$  of the  $i$ -th response is fixed after time  $t_i$  so  $Y_i = Y_i(t_i)$ . Additionally, it is supposed that model (1) operates in the so-called external environment, which has final state space  $S = \{s_j, j = 1, \dots, m\}$ . In our case  $m=2$ . For each state parameters  $\beta$  are different:  $\beta_1 = (\beta_{1,1}, \beta_{2,1}, \dots, \beta_{k,1})^T$  and  $\beta_2 = (\beta_{1,2}, \beta_{2,2}, \dots, \beta_{k,2})^T$ , where  $k$  is a number of factors. Let  $\vec{t} = (t_{i,1}, t_{i,2})$  be  $1 \times 2$  vector, for that component  $t_{i,j}$  means a sojourn time for response  $Y_i$  in the state  $j$  (thus,  $t_i = t_{i,1} + t_{i,2}$ ). Then (taking into account properties of the normal distribution) we can rewrite formula (1):

$$Y_i(t_i) = x_i(\beta_1 t_{i,1} + \beta_2 t_{i,2}) + Z_i \sqrt{t}. \quad (2)$$

Further let us define the model in matrix notation

$$Y(t) = (Y_1(t_1), \dots, Y_n(t_n))^T = \begin{pmatrix} \vec{t}_1 \otimes x_1 \\ \vec{t}_2 \otimes x_2 \\ \dots \\ \vec{t}_n \otimes x_n \end{pmatrix} \text{vec } \beta + \text{diag}(\sqrt{t_1}, \sqrt{t_2}, \dots, \sqrt{t_n})Z,$$

where  $Y_i(t)$  are scale responses which are time-additive ( $Y_i(0) = 0$ ),  $n$  is the number of observations, the  $1 \times 2$  vector  $\vec{t}_i = (t_{i,1}, t_{i,2})$ , the  $n \times 2$  matrix  $T = (\vec{t}_1^T, \dots, \vec{t}_n^T)^T$ ,  $\otimes$  is Kronecker product,  $(x_{i,1}, x_{i,2}, \dots, x_{i,k})$  is  $1 \times k$  vector, the  $k \times 2$  matrix  $\beta = (\beta_1, \beta_2) = (\beta_{v,j})$  of unknown parameters,  $\text{vec}$  operator  $\text{vec}(A)$  of matrix  $A$ , the  $n$ -dimensional diagonal matrix  $\text{diag}(v)$  with the vector  $v$  on the main diagonal,  $Z = (Z_i)$  is the  $n \times 1$  vector, where  $Z_i(t)$  is Brown motion scale disturbance ( $Z_i$  are independently, identically normally distributed with mean zero and constant variance  $\sigma^2$ ).

This is the case of the generalized linear regression model. The whole trajectory of the environment  $J(\cdot)$  is unknown and the estimated conditional average sojourn time is used instead of unknown sojourn times  $T_{i,j}$  in the state  $s_j$ . The solution of a usual system of differential equations for Markov chain is represented by the matrix exponent. And if all the eigenvalues of matrix are different, then the simpler solution is obtained for the conditional average sojourn time calculation. Further unknown parameters  $\beta$  are estimated as follows:

$$\text{vec } \tilde{\beta} = \left( \sum_{i=1}^n \frac{1}{t_i} (\vec{t}_i^T \vec{t}_i) \otimes (x_i^T x_i) \right)^{-1} \cdot \begin{pmatrix} t_1^{-1} \vec{t}_1 \otimes x_1 \\ t_2^{-1} \vec{t}_2 \otimes x_2 \\ \dots \\ t_n^{-1} \vec{t}_n \otimes x_n \end{pmatrix}^T Y$$

All necessary formulas for a calculation of the conditional average sojourn time that allows to get the needed estimates and, also full description of the Markov-modulated linear regression model are provided in previous researches [1, 2].

### III. EXTERNAL ENVIRONMENT: DATA DESCRIPTION

Weather conditions in the city of Ventspils were chosen as an external environment. It was assumed that there are two states of the external environment: “No precipitation” or “dry” weather conditions and “Precipitation” or “wet” (alternating states). Division was made subjectively, but based on research objectives, namely, all weather conditions that worsen visibility and may influence different transport indicators (punctuality in this case) belong to the second group (“wet”), and, accordingly, the rest belongs to the first group (“dry”). The included data about weather conditions was obtained from the Latvian Environment, Geology and Meteorology Centre (LEGMC) database. General information regarding this resource is paraphrased from the LEGMC homepage [www.meteo.lv](http://www.meteo.lv), while details related to the included weather data are presented below.

Meteorological observations are carried out by LEGMC at 33 observation stations (1 station per 1500km<sup>2</sup>), which are stationary and located over the territory of Latvia. Stations location is optimal to provide a sufficiently detailed description of Latvia weather conditions and climate. For this research data from the station named “Ventspils” which is located on the west coast of Latvia in the Ventspils city was used. All available data can be downloaded from LEGMC website in an excel format.

Monthly data about “Past weather conditions 1” was processed from 1986 to 2017 inclusively (except years 2010, 2011, 2012 – data was not available). Data was selected by parameter “Past weather conditions 1”, which contains codes from 0 to 9, with measurement points every 3 hours. Due to LEGMC explanations the results of visual observations are recorded in coded form. The codes are represented in Table 1.

Table 1. Codes of selected parameter

| Parameter: past weather conditions 1 and 2 |  |
|--|--|
| <b>0</b>                                   | The amount of clouds is less than 5 points between observation boundaries, clear   |
| <b>1</b>                                   | The amount of clouds has changed from $< 5$ to $> = 5$ points between observation boundaries   |
| <b>2</b>                                   | The amount of clouds covers $> 5$ points between observation boundaries  |
| <b>3</b>                                   | All kinds of drift storms (snow drifting close to the ground with or without snow falling) between the observation boundaries. Dust storms or sandstorms between observation boundaries. |
| <b>4</b>                                   | Mist or ice fog between observation boundaries. Smog between observation boundaries. Visibility $< 1$ km   |
| <b>5</b>                                   | Drizzle between observation boundaries   |
| <b>6</b>                                   | Rain between observation boundaries  |
| <b>7</b>                                   | Snow, snow pellets, needle ice or ice pellets, rain with snow between observation boundaries   |
| <b>8</b>                                   | Heavy precipitation (heavy snowfall, downpour, snow or ice grains, hail) between observation boundaries  |
| <b>9</b>                                   | Thunderstorm with or without precipitation between observation boundaries  |

All weather condition codes were divided into two groups mentioned above: “No precipitation” or “dry” weather conditions: codes from 0 to 2, and “Precipitation” or “wet”: codes from 3 to 9.

The Visual Basic code was written to divide the available codes (from 0 – 9) into two groups and to calculate the duration of a sojourn time in each state (multiplying the number of consecutive values in each group by 3 (in hours)). Thus, obtaining data to estimate the distribution of the sojourn time in each state.

It is natural to assume that different seasons will have distinctive characteristics of the transition intensities from state to state. At the first stage it was decided to look into the autumn and winter months: September, October, November, December, January and February. Descriptive and inferential analysis was carried out by means of statistical software package Statistica 12.0. Tables 2.1 and 2.2 represent descriptive statistics of the sojourn time in each state for each autumn and winter month accordingly.

Table 2.1. Descriptive statistics of the sojourn time in each state for each autumn month

| Variable  | Valid N | Mean     | Sum   | Max | Std.Dev. |
|-----------|---------|----------|-------|-----|----------|
| Sept/NoPr | 507     | 29.81065 | 15114 | 372 | 43.38    |
| Sept/Prec | 492     | 11.76220 | 5787  | 99  | 10.30010 |
| Oct/NoPr  | 557     | 23.39677 | 13032 | 237 | 33.66277 |
| Oct/Prec  | 548     | 14.32117 | 7848  | 93  | 11.60828 |
| Nov/NoPr  | 628     | 16.98726 | 10668 | 141 | 21.06123 |
| Nov/Prec  | 617     | 16.55105 | 10212 | 105 | 14.76343 |

Table 2.2. Descriptive statistics of the sojourn time in each state for each winter month

| Variable | Valid N | Mean     | Sum   | Max | Std.Dev. |
|----------|---------|----------|-------|-----|----------|
| Dec/NoPr | 644     | 18.02329 | 11607 | 162 | 21.00399 |
| Dec/Prec | 635     | 15.65669 | 9942  | 99  | 14.48611 |
| Jan/NoPr | 615     | 19.25366 | 11841 | 153 | 23.55728 |
| Jan/Prec | 617     | 15.77796 | 9735  | 153 | 14.91429 |
| Feb/NoPr | 545     | 21.45138 | 11691 | 210 | 26.80865 |
| Feb/Prec | 541     | 14.72274 | 7965  | 105 | 13.04184 |

It is evident from the Table 2 and the Fig.1 that the average sojourn time in the "dry" state decreases from the first autumn month to the last, and vice versa the average sojourn time in the "dry" state increases from the first winter month to the last, which seems to be natural and therefore can serve as a validation of the conceptual model.

Other issues to be addressed: is it possible to combine all three months to describe the behaviour of the external environment? Are transition intensities the same for all autumn months? Answering these questions requires that homogeneity analysis should be carried out. The null hypothesis (in words) in general can be stated as follows: the sojourn time in particular state ("dry" or "wet") is identical for Month1 and Month2. More formally,  $H_0: F_1(x_1, \dots, x_n) = F_2(y_1, \dots, y_n)$ . Two nonparametrical tests were used for homogeneity analysis at the significance level of 0.01: Mann-Whitney U test and Kolmogorov-Smirnov two-sample test. The results for different combinations of parameters are shown in the Table 3.

With the chosen significance level (0.01), just two hypotheses for the autumn months were accepted. With increasing significance level (for example, 0.05), all hypotheses would be rejected. Based on the obtained results, it

was decided to carry out an experiment on the combined sample with the weather data for October and November.

As for the winter months, the data showed that the difference in the average sojourn time of weather conditions in one of the states (precipitation or no precipitation) can be considered insignificant for the three winter months. Thus, the data for the winter months can be considered homogeneous and can be combined for further analysis.

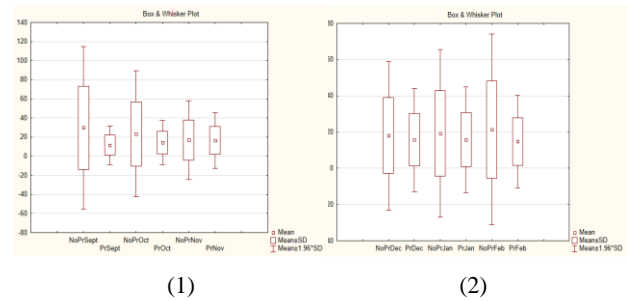


Fig.1. Box&Whisker plot of the sojourn time in each state for each autumn (1) and winter (2) month (See much better resolution of Fig 1 (1), (2) in the Appendix)

Table 3. Homogeneity testing results

| H <sub>0</sub> : State, Month1 vs Month2 | Mann-Whitney Test Z/ p-value | Kolmogorov-Smirnov test/ p-value | Interpretation of the results |
|--|------------------------------|----------------------------------|-------------------------------|
| “No prec”, Sept vs Oct                   | 2.981851/0.002865            | 0.096/p < .025                   | Reject H <sub>0</sub>         |
| “Precipitation” Sept vs Oct              | -3.80089/0.000144            | -0.116/p < .005                  | Reject H <sub>0</sub>         |
| “No prec”, Oct vs Nov                    | 2.196350/0.028068            | 0.071567/p < .10                 | Accept H <sub>0</sub>         |
| “Precipitation” Oct vs Nov               | -2.21623/0.026676            | -0.058231/p > .10                | Accept H <sub>0</sub>         |
| “No prec”, Sept vs Nov                   | 5.305069/0.000000            | 0.159151/p < .001                | Reject H <sub>0</sub>         |
| “Precipitation” Sept vs Nov              | -6.06640/0.0000              | -0.174144/p < .001               | Reject H <sub>0</sub>         |
| “No prec”, Dec vs Jan                    | -0.676/0.4989                | -0.029/p > .10                   | Accept H <sub>0</sub>         |
| “Precipitation” Dec vs Jan               | -0.469/0.638373              | -0.027/p > .10                   | Accept H <sub>0</sub>         |
| “No prec”, Jan vs Feb                    | -1.38246/0.166833            | -0.051421/p > .10                | Accept H <sub>0</sub>         |
| “Precipitation” Jan vs Feb               | 1.408605/0.158953            | 0.049860/p > .10                 | Accept H <sub>0</sub>         |
| “No prec”, Dec vs Jan                    | -2.04293/0.041060            | -0.055171/p > .10                | Accept H <sub>0</sub>         |
| “Precipitation” Dec vs Feb               | 0.897714/0.369339            | 0.045789/p > .10                 | Accept H <sub>0</sub>         |

Two criteria were selected to check markovian properties of the described external environment: distribution of the sojourn time in each state (which is supposed to be exponential) and independence of the observations' pairs (memoryless property).

Visualization of the sojourn time in each state for each autumn month by means of histogram (Fig.2) showed satisfactory results and suggests that the distribution of the sojourn time is indeed exponential. Since it was decided to

consider the sojourn time in each state to be homogeneous for all winter months, the histograms of the distributions were built on the combined data (Fig.3), which also implies the presence of an exponential distribution. It is necessary to test hypothesis about the sojourn time distribution in each state. The null hypothesis (in words) in general can be stated as follows: the sojourn time in particular state (“No precipitation” or “Precipitation”) has an exponential distribution. More formally,  $H_0: F_{emp}(x_1, \dots, x_n) = F_{exp}(x_1, \dots, x_n)$ . Two nonparametrical tests were used in distribution fitting procedure: Chi-square test and Kolmogorov-Smirnov test. The results for different combinations of parameters at the significance level 0.01 are shown in the Table 4 for autumn months and in the Table 5 for combined winter months.

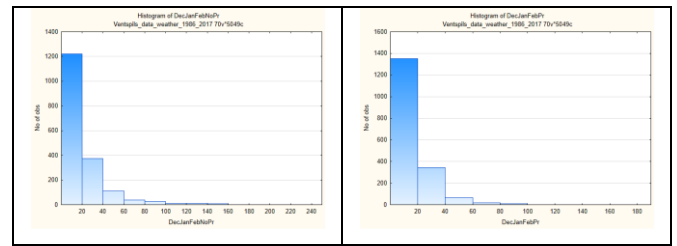


Fig.3. Histograms of the sojourn time in each state for all winter months combined (See much better resolution of Fig 3 in the Appendix)

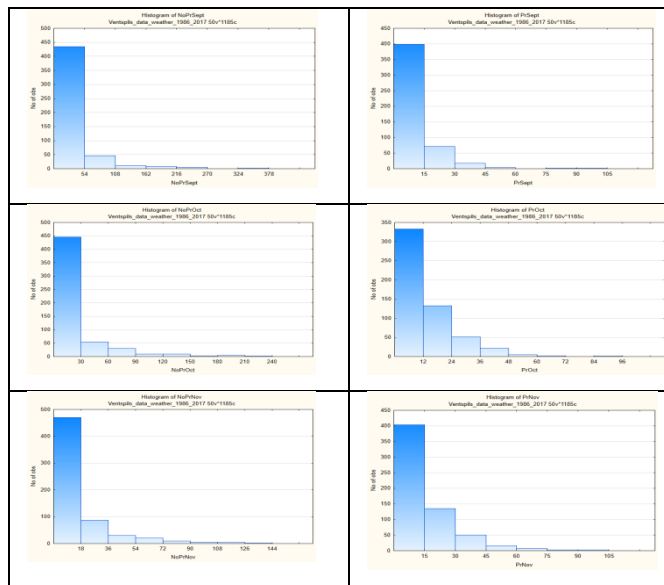


Fig.2. Histograms of the sojourn time in each state for each autumn month (See much better resolution of Fig 2 in the Appendix)

Table 4. Distribution fitting results (1)

| H <sub>0</sub> : State, Month | Chi-Square test / p-value | Kolmogorov-Smirnov test/ p-value | Interpretation of the results |
|-------------------------------|---------------------------|----------------------------------|-------------------------------|
| “No prec”, September          | 14.65/ 0.00013            | 0.16369/ p < 0.01                | Reject H <sub>0</sub>         |
| “Precipitation” September     | 6.90211/ 0.14115          | 0.22513/ p < 0.01                | Accept H <sub>0</sub>         |
| “No prec”, October            | 50.9394/ 0.0000           | 0.17619/ p < 0.01                | Reject H <sub>0</sub>         |
| “Precipitation” October       | 10.58989/ 0.10191         | 0.18899/ p < 0.01                | Accept H <sub>0</sub>         |
| “No prec”, November           | 36.3762/ 0.00000          | 0.16189/ p < 0.01                | Reject H <sub>0</sub>         |
| “Precipitation” November      | 5.72381/ 0.33403          | 0.16793/ p < 0.01                | Accept H <sub>0</sub>         |
| “No prec”, Oct+Nov            | 104.582/ 0.00000          | 0.15746/ p < 0.01                | Reject H <sub>0</sub>         |
| “Precipitation” Oct+Nov       | 10.0166/ 0.0747           | 0.17595/ p < 0.01                | Accept H <sub>0</sub>         |

Table 5. Distribution fitting results (2)

| H <sub>0</sub> : State, Month | Chi-Square test / p-value | Kolmogorov-Smirnov test/ p-value | Interpretation of the results |
|-------------------------------|---------------------------|----------------------------------|-------------------------------|
| “No prec”, December           | 81.85981/ 0.0000          | 0.14274/ p < 0.01                | Reject H <sub>0</sub>         |
| January                       |                           |                                  |                               |
| February                      |                           |                                  |                               |
| “Precipitation” December      | 27.84398/ 0.00023         | 0.19858/ p < 0.01                | Reject H <sub>0</sub>         |
| January                       |                           |                                  |                               |
| February                      |                           |                                  |                               |

Distribution fitting procedure showed that for each considered autumn month the null hypothesis about exponentiality of the sojourn time in the state “Precipitation” couldn’t be rejected for the given sample and at the chosen significance level. However, the same null hypothesis for the state “No precipitation” was rejected. The same hypothesis rejection situation was revealed for the sample with combined winter months. Even though the results obtained were partially negative, it was decided that the data is relevant for describing the external environment in the context of this experiment. Moreover, it does not seem to be a problem, since author’s current studies prove using an approximation of arbitrary nonnegative density by a convolution of exponential densities.

Correlation analysis was made only for autumn months and indicated that there is no linear dependence between the observation pairs "No precipitation – Precipitation" for each month. The results are presented in Table 6.

Table 6. Correlation analysis results

| Variable | Correlations (Casewise deletion of missing data) |          |                  |                  |
|----------|--|----------|------------------|------------------|
|          | Means  | Std.Dev. | NoPrSept         | PrSept           |
| N=492    |  |          |                  |                  |
| NoPrSept | 29.85366   | 43.45491 | 1.000000         | <b>0.057174</b>  |
| PrSept   | 11.76220   | 10.30010 | <b>0.057174</b>  | 1.000000         |
| N=548    |  |          |                  |                  |
| NoPrOct  | 23.28285   | 33.16281 | 1.000000         | <b>0.098189</b>  |
| PrOct    | 14.32117   | 11.60828 | <b>0.098189</b>  | 1.000000         |
| N=617    |  |          |                  |                  |
| NoPrNov  | 17.12966   | 21.20255 | 1.000000         | <b>-0.033114</b> |
| PrNov    | 16.55105   | 14.76343 | <b>-0.033114</b> | 1.000000         |

Data on means and standard deviations are partly different from the data presented in Table 2 (3 out of 6 values for each indicator differ insignificantly). It stems from deliberate exclusion of the data for which no corresponding pair was

found in the analysis of the relations between the two variables.

In general, the test results of the external environment model for Markovian property are subject to different interpretations but has been considered as admissible ones for the experimental purposes.

#### IV. COACHES' DELAY TIME: DATA DESCRIPTION

Ventspils-Riga route was selected for the analysis of coaches' delay time. Depending on the day of the week, from Monday to Sunday, there are 16, 15, 14, 14, 14, 18 and 13 scheduled runs, correspondingly. The scheduled duration of the run is also a variable and ranges from 180 to 255 minutes. The RICT management provided data on scheduled and actual departure and arrival time of the coaches as well as the record date and capacity of a coach. This data covers the period from 2012 to 2017. For example, the data for 2012 contains 5414 records. The above-mentioned statistical data was made suitable for the analysis. Since in the Markov-modulated linear regression model the dependent variable is time-additive, the delays were summed for each day of the week, a total of 365 observations were obtained for 2013, 2014, 2015 and 2017 years, and 366 observations for 2012 and 2016 leap years, respectively. Since at the first stage of the analysis only autumn and winter months are considered, the sample size has naturally decreased. The given sample can also be analysed as time series with various patterns characteristic of this type of observation, but this task goes beyond the scope of the current study.

One of the main principles of RICT management is the provision of quality services. Punctuality is an essential component of the quality system. The bulk of the delays is within acceptable limits, however in any system unforeseen circumstances may arise and cause schedule shifts (and, as a result, delays) of coaches. Thus, according to the results of a survey of coaches' drivers about the factors having negative impact on adherence to a timetable (RICT, internal procedure D07, 2017), 70% indicated traffic jams in Riga, 16% - weather conditions, 3% - technical condition of the coach, 8% - coach route timetable, and 3% indicated other reasons.

Table 7 presents the average total delays for each day of the week (according to the 6 years sample). According to the results for all autumn and winter months, the least successful day from the point of view of punctuality is Friday, and the most successful ones are Saturday and Sunday, which can also be regarded as confirmation of the data validity. Firstly, since many studies show that the traffic intensity is higher on Fridays, for example, [5, 6], and, secondly, due to the general human experience that confirms this fact.

Table 7. Average delay time of coaches distributed by days

|               | Average delay time, minutes |              |              |              |              |              |              |
|---------------|-----------------------------|--------------|--------------|--------------|--------------|--------------|--------------|
|               | Mon                         | Tue          | Wed          | Thu          | Fri          | Sat          | Sun          |
| Sept<br>N=180 | 44.3<br>N=26                | 36.8<br>N=26 | 35.4<br>N=25 | 34.5<br>N=25 | 78.6<br>N=26 | 14.6<br>N=26 | 18.0<br>N=26 |
| Oct<br>N=186  | 37.1<br>N=27                | 21.8<br>N=27 | 21.4<br>N=27 | 22.2<br>N=27 | 52.2<br>N=26 | 12.6<br>N=26 | 10.8<br>N=26 |
| Nov<br>N=180  | 19.4<br>N=25                | 24.1<br>N=25 | 22.4<br>N=26 | 29.6<br>N=26 | 47.8<br>N=26 | 8.8<br>N=26  | 10.0<br>N=26 |
| Dec<br>N=186  | 26.9<br>N=27                | 31.0<br>N=27 | 22.0<br>N=26 | 21.0<br>N=26 | 40.8<br>N=26 | 6.5<br>N=27  | 10.2<br>N=27 |

|              |              |              |              |              |              |             |             |
|--------------|--------------|--------------|--------------|--------------|--------------|-------------|-------------|
| Jan<br>N=186 | 12.6<br>N=26 | 13.3<br>N=27 | 15.3<br>N=26 | 24.4<br>N=27 | 29.6<br>N=27 | 7.9<br>N=26 | 6.0<br>N=27 |
| Feb<br>N=170 | 10.4<br>N=25 | 7.0<br>N=24  | 9.4<br>N=25  | 8.7<br>N=24  | 20.5<br>N=24 | 5.4<br>N=24 | 6.8<br>N=24 |

Also, from the table 7 it is clear that in the winter months punctuality increases, but situation with Friday remains the same.

#### V. MODELLING RESULTS

Modelling was made by means of Mathcad 14.0 and Excel® software.

Four main experiments were carried out. The structure of the input data is the same for all experiments, only the values differ. To apply the model, it is necessary to have the following initial data: the matrix of the state transition intensities ( $\lambda$ ), the matrix with the regressors' values ( $X$ ), the vector of observation durations ( $\tau$ ), the vector with the external environment initial states ( $I$ ), and the vector of dependent variable values ( $Y$ ).

A random environment has two states ( $m=2$ ). The transition intensities from state  $i$  to state  $j$  are calculated as reciprocals to the sojourn time. The days of the week serve as regressors. The number of regressors is seven: six dummy variables, and a constant term. Days of the week are represented as dummy variables, "Monday" serves as a key variable.

##### A. Experiment 1

Data on coaches' delay times for the month of September from 2013 to 2017 was used for experiment 1. (Weather data for 2012 is not available. Consequently, there is no possibility to plot a vector with external environment initial states; therefore, data on delays for 2012 is excluded from consideration).

Transition rates from state  $i$  to state  $j$  are given by the transition matrix:

$$\lambda := \begin{pmatrix} 0 & 0.033545058 \\ 0.085018144 & 0 \end{pmatrix}$$

Stationary probabilities of states are as follows:

$$\pi = (0.717 \quad 0.283)^T$$

Dimensions of given matrices and vectors:

$$X_{150 \times 7}, Y_{150 \times 1}, \tau_{150 \times 1}, I_{150 \times 1}$$

We begin with the estimates for the simple linear regression (ordinary weighted least squares) with 7 regressors:  $vec(\tilde{\beta}) = (44.318 \ -7.136 \ -6.556 \ -6.747 \ 30.182 \ -28.604 \ -26.985)$ , with noticeably large residual sum of squares  $RSS = 242700$  and determination coefficient  $R\text{-squared} = 0.17$ .

Further we use supposed approach and get estimations with respect the external environment. Since we have two states and seven independent variables, the number of unknown parameters  $\beta$  equals to 14.

$vec(\tilde{\beta}) = (4.217 \ -1.254 \ -0.663 \ -0.854 \ 2.124 \ -3.498 \ -2.587 \ -0.055 \ 1.92 \ -1.093 \ 0.042 \ 1.642 \ 2.242 \ -0.024)$  with smaller  $RSS = 230700$  and higher determination coefficient  $R\text{-squared} = 0.211$ .  $RSME = 39.346$ .

Compering to Monday (which is key variable) for all other days (except Friday) according to coefficients which have negative sign, delays are smaller (for the first state "No

precipitation”), which seems adequate according to Table 6. For the second state “Precipitation”, the coefficients appear with more random signs (what is less explainable).

### B. Experiment 2

Since explicit validation set is not available 6-fold cross-validation technique was used for assessing accuracy of model prediction. Each validation set consists of  $n = 25$  observations.

Transition rates from state  $i$  to state  $j$  and stationary probabilities of states are the same as in experiment 1.

Dimensions of given matrices and vectors:  $X_{125 \times 7}$ ,  $Y_{125 \times 1}$ ,  $\tau_{125 \times 1}$ ,  $I_{125 \times 1}$ .

Table 8 contains the results of 3 iterations of model estimation.

Table 8. Estimations of unknown parameters, 3 iterations

| Parameter            | Iteration 1   | Iteration 2   | Iteration 3   |
|----------------------|---------------|---------------|---------------|
| $\tilde{\beta}_{00}$ | 4.613         | 5.419         | 4.144         |
| $\tilde{\beta}_{01}$ | -1.033        | -1.917        | -1.069        |
| $\tilde{\beta}_{02}$ | -0.434        | -1.472        | -0.535        |
| $\tilde{\beta}_{03}$ | -1.187        | -1.267        | -0.496        |
| $\tilde{\beta}_{04}$ | 1.874         | 1.32          | 2.439         |
| $\tilde{\beta}_{05}$ | -3.952        | -4.537        | -3.383        |
| $\tilde{\beta}_{06}$ | -2.605        | -3.618        | -2.227        |
| $\tilde{\beta}_{10}$ | 0.144         | -0.802        | -0.245        |
| $\tilde{\beta}_{11}$ | 1.601         | 2.513         | 2.969         |
| $\tilde{\beta}_{12}$ | -1.828        | -0.578        | -0.935        |
| $\tilde{\beta}_{13}$ | 0.792         | 0.324         | 0.065         |
| $\tilde{\beta}_{14}$ | 2.389         | 2.154         | 2.127         |
| $\tilde{\beta}_{15}$ | 2.646         | 2.892         | 2.407         |
| $\tilde{\beta}_{16}$ | -1.252        | 0.628         | 7.624e-3      |
| RSS                  | 213900        | 205800        | 218400        |
| RSS*                 | 20230         | 30410         | 13640         |
| RMSE                 | <b>41.366</b> | <b>40.572</b> | <b>41.802</b> |
| RMSE*                | <b>28.448</b> | <b>34.876</b> | <b>23.354</b> |

\* - out-of-sample (testing sample)

All iterations showed not so high out-of-sample prediction power of the model, but RMSE was smaller for all cases.

### C. Experiment 3

Data on coaches’ delay times for the month of October from 2013 to 2017 was used for experiment 2.

Transition rates from state  $i$  to state  $j$  are given by the transition matrix:

$$\lambda := \begin{pmatrix} 0 & 0.0427 \\ 0.0698 & 0 \end{pmatrix}.$$

Stationary probabilities of states are as follows:

$$\pi = (0.62 \quad 0.38)^T.$$

105 observations were used as a training set and the rest of 50 observations as a validation set.

Dimensions of given matrices and vectors:

$X_{105 \times 7}$ ,  $Y_{105 \times 1}$ ,  $\tau_{105 \times 1}$ ,  $I_{105 \times 1}$ .

In the Table 9 comparison of the expectation of the responses ( $E(Y)$ ) and actual responses ( $Y$ ) for validation set is shown.

Table 9. Comparison of actual and expected responses Y

| $n$ | 0 | 1  | 2  | 3  | 4 | 5 | 6 | ... |
|-----|---|----|----|----|---|---|---|-----|
| $Y$ | 0 | 12 | 10 | 40 | 3 | 0 | 5 | ... |

|        |      |      |      |      |      |      |      |     |
|--------|------|------|------|------|------|------|------|-----|
| $E(Y)$ | 21.7 | 16.2 | 28   | 54.3 | 13.5 | 13.4 | 48   | ... |
| $n$    | 28   | 29   | 30   | 31   | 32   | 33   | 34   | ... |
| $Y$    | 5    | 3    | 15   | 44   | 22   | 54   | 0    | ... |
| $E(Y)$ | 16.2 | 28   | 33.8 | 11.5 | 54.3 | 15.7 | 13.4 | ... |

Table 10 contains estimated model quality criteria such as RSS and RMSE. Results within the validation set showed satisfactory results.

Table 10. Model quality criteria

| Type of set | RSS    | RMSE   |
|-------------|--------|--------|
| Training    | 114500 | 33.017 |
| Validation  | 20660  | 20.328 |

### D. Experiment 4

Given experiment combines data of two autumn months: October and November, thus expanding the sample twice.

Transition rates from state  $i$  to state  $j$  are given by the transition matrix:

$$\lambda := \begin{pmatrix} 0 & 0.05 \\ 0.064507198 & 0 \end{pmatrix}.$$

Stationary probabilities of states are as follows:

$$\pi = (0.563 \quad 0.437)^T.$$

Sample size is equal to 305. Estimation gives the following results: RSS = 241400, RMSE = 28.18 and R-squared = 0.173. If we compare with all the experiments, the last model showed the best results based on RMSE criterion.

## VI. CONCLUSION

This paper considers implementation of Markov-modulated linear regression model into practice. Actual data on coaches’ trip times is provided by the Riga International Coach Terminal. Data on weather conditions provided by the Latvian Environment, Geology and Meteorology Centre and is free downloaded from database on the website. Preliminary data preparation was carried out both for external environment description and regression model development itself.

Data on weather conditions in the Ventspils city is used for the external environment description: two states are assumed: “no precipitation” and “precipitation”. The average sojourn time in the state “No precipitation” decreases from the first autumn month to the last and vice versa increases from the first winter month to the last. The model of the external environment is tested for the markovian properties. Two criteria were selected to check markovian properties of the described external environment: distribution of the sojourn time in each state (which is supposed to be exponential) and independence of the observations’ pairs. Despite the fact that the results obtained were partially negative, it was decided that the data is relevant for describing the external environment in the context of this experiment. Actual data on coaches’ trip times was analysed by means of descriptive statistics. The least successful day from the point of view of punctuality is Friday, and the most successful ones are Saturday and Sunday, which seems quite natural.

The application of Markov-modulated linear regression model on this sample did show quite adequate results. The low

accuracy of the prediction is not related to the incorrectness of the proposed model, rather it is due to the low quality of the model for describing the delays of coaches: supposedly the day of the week is not the only factor determining the size of the delay. One more reason relates to the quality and amount of data: matrices  $\mathbf{X}$  and vector  $\mathbf{I}$  are sparse matrix and vector, also vector  $\boldsymbol{\tau}$  contains mostly repeating elements, that could cause unreliable results. The following tasks are requested for further investigation:

- Try different k-folds within cross-validation technique obtaining more reliable results.
- To analyse the remaining months of the year (spring and summer).
- To consider inclusion of other factors in the model (as independent variables).

Markov-modulated linear regression model has recommended itself positively and requires further development.

### **Acknowledgment I**

This work was financially supported by the specific support objective activity 1.1.1.2. "Post-doctoral Research Aid" (Project id. N. 1.1.1.2/16/I/001) of the Republic of Latvia, funded by the European Regional Development Fund. NADEZDA SPIRIDOVSKA RESEARCH PROJECT NO. 1.1.1.2/VI AA/1/16/075 "NON-TRADITIONAL REGRESSION MODELS IN TRANSPORT MODELLING"

### **Acknowledgment II**

The author appreciates the management of the Riga International Coach Terminal for provided data.

### **Acknowledgment III**

The author would like to thank the Associate Editor for helping her in technical details and the Reviewers

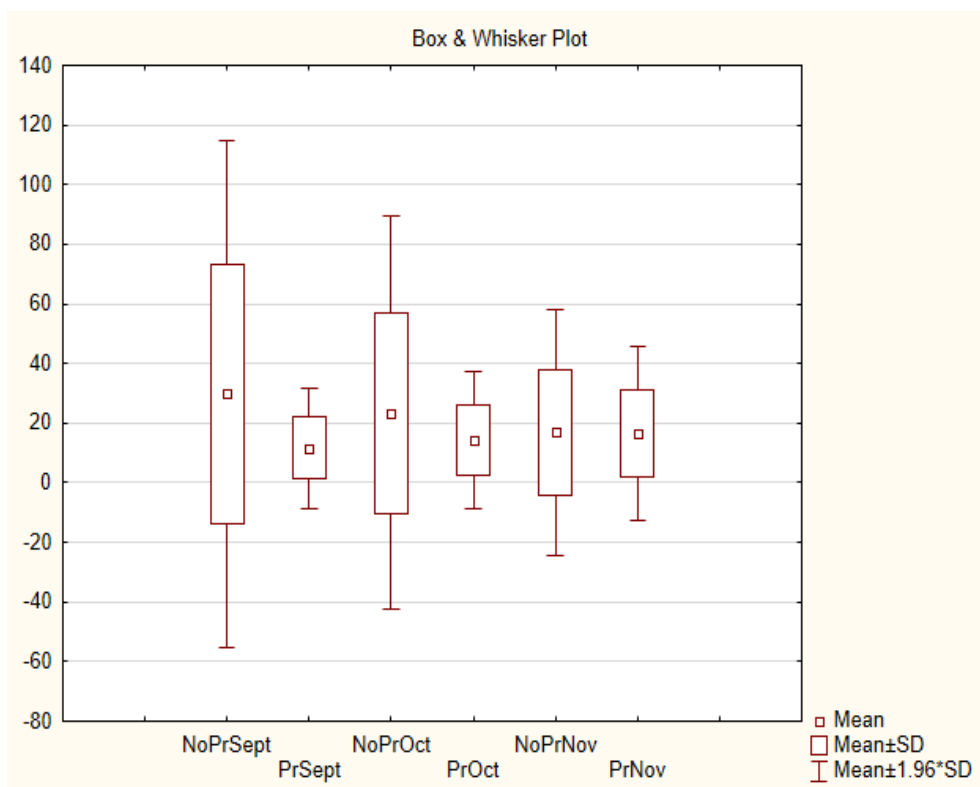
### **References**

- [1] Andronov, A., Spiridovska, N. Markov-Modulated Linear Regression. In proceedings' book: International conference on Statistical Models and Methods for Reliability and Survival Analysis and Their Validation (S2MRSA), 2012, pp.24–28. Bordeaux, France
- [2] Andronov A. Parameter statistical estimates of Markov-modulated linear regression, in: Statistical Methods of Parameter Estimation and Hypothesis Testing 24, Perm State University, Perm, Russia, 2012, pp. 163–180. (Russian).
- [3] Spiridovska N. Markov-Modulated Linear Regression: Tasks and Challenges in Transport Modelling. In: Kabashkin I., Yatskiv I., Prentkovskis O. (eds) Reliability and Statistics in Transportation and Communication. RelStat 2017. Lecture Notes in Networks and Systems, vol 36. Springer, Cham, 2018.
- [4] I. Yatskiv. (Jackiva), E. Budilovich. (Budiloviča), and V. Gromule, "Accessibility to Riga Public Transport Services for Transit Passengers," Procedia Eng., vol. 187, pp. 82–88, 2017.

- [5] Pacheco A., Tang L.C., Prabhu N.U. Markov-Modulated Processes & Semiregenerative Phenomena. New Jersey – London; World Scientific, 2009.
- [6] E. Steiger, B. Resch, J. P. de Albuquerque, and A. Zipf, "Mining and correlating traffic events from human sensor observations with official transport data using self-organizing-maps," Transp. Res. Part C Emerg. Technol., vol. 73, pp. 91–104, Dec. 2016.
- [7] N. Earl, I. Simmonds, and N. Tapper, "Weekly cycles in peak time temperatures and urban heat island intensity," Environ. Res. Lett., vol. 11, pp. 074003-, Jul. 2016.
- [8] M. A. Rotondi, "To Ski or Not to Ski: Estimating Transition Matrices to Predict Tomorrow's Snowfall Using Real Data," J. Stat. Educ., vol. 18, no. 3, 2010.
- [9] P. Jordan and P. Talkner, "A seasonal Markov chain model for the weather in the central Alps," Tellus Dyn. Meteorol. Oceanogr., vol. 52, no. 4, pp. 455–469, 2000.
- [10] Kijima M. Markov Processes for Stochastic Modeling. London: Chapman & Hall, 1997.

### **APPENDIX I**

In the Appendix, we present the Figures 1 and 2 with better resolution. The Appendix has been added by the Publisher to facilitate the readers of the article.



(1)

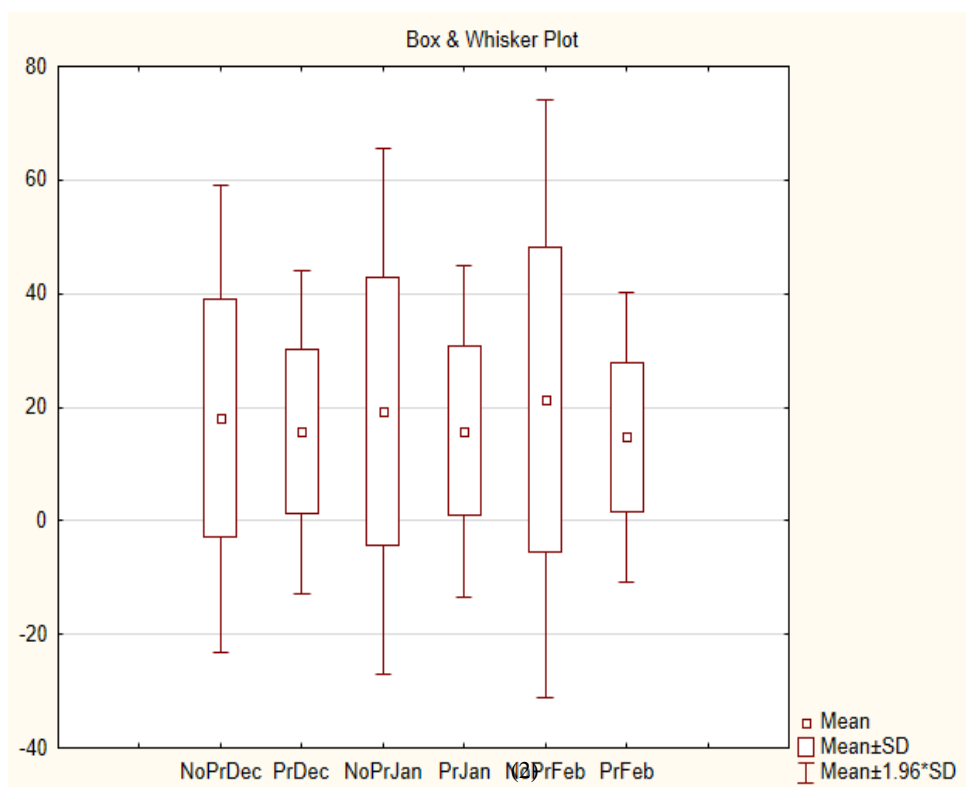
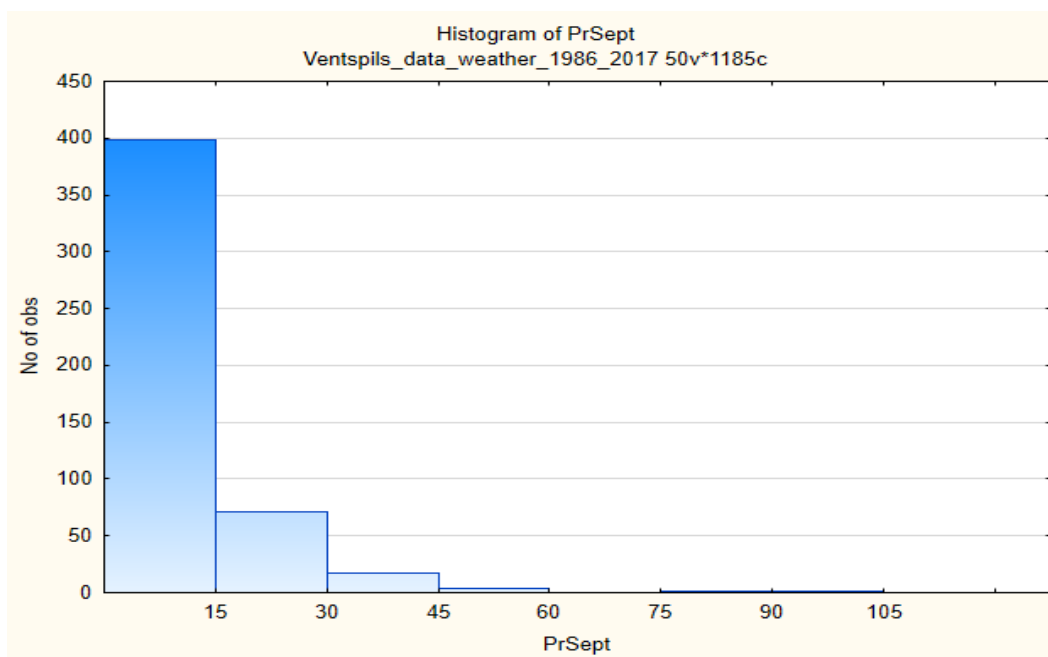
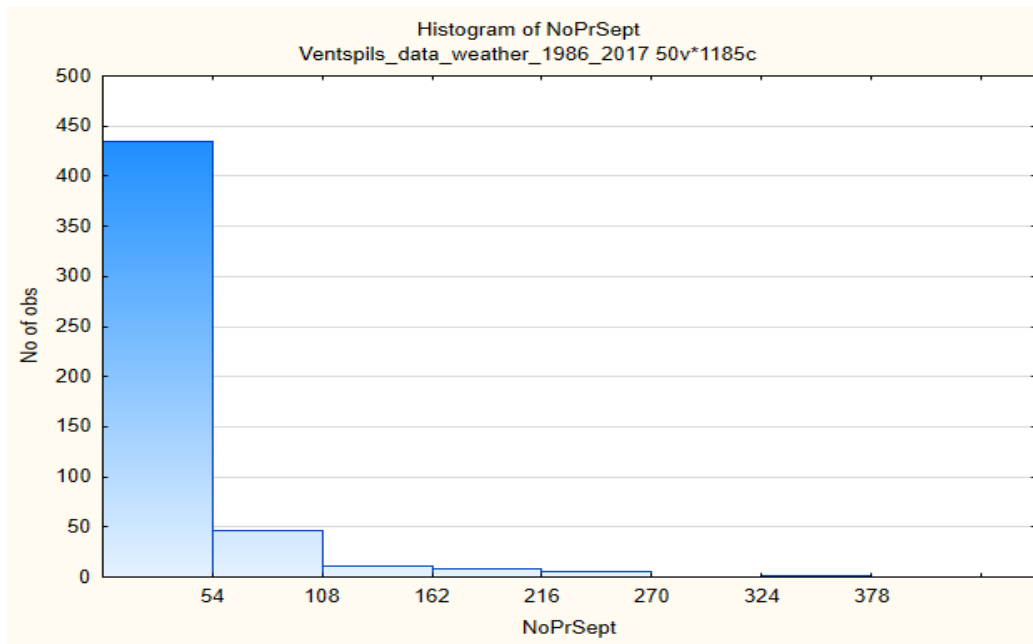
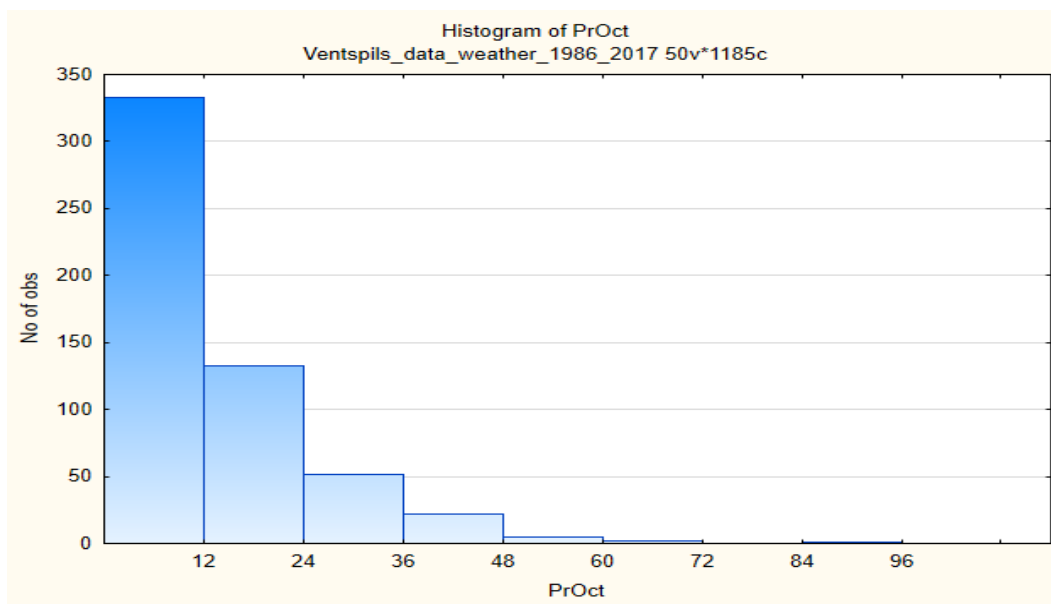
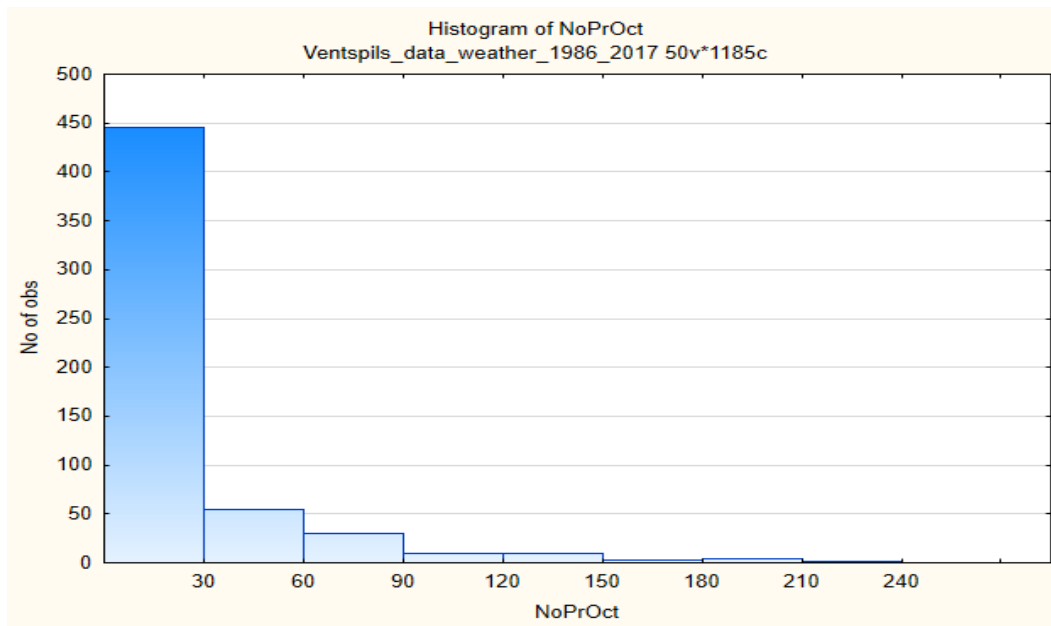


Fig.1. Box&Whisker plot of the sojourn time in each state for each autumn (1) and winter (2) month







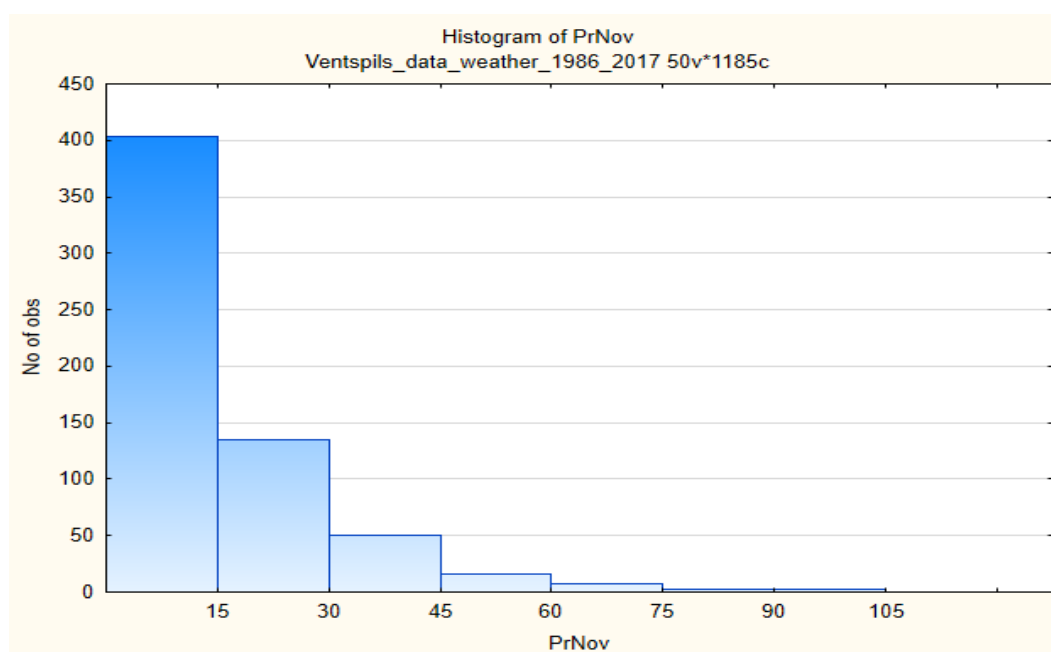
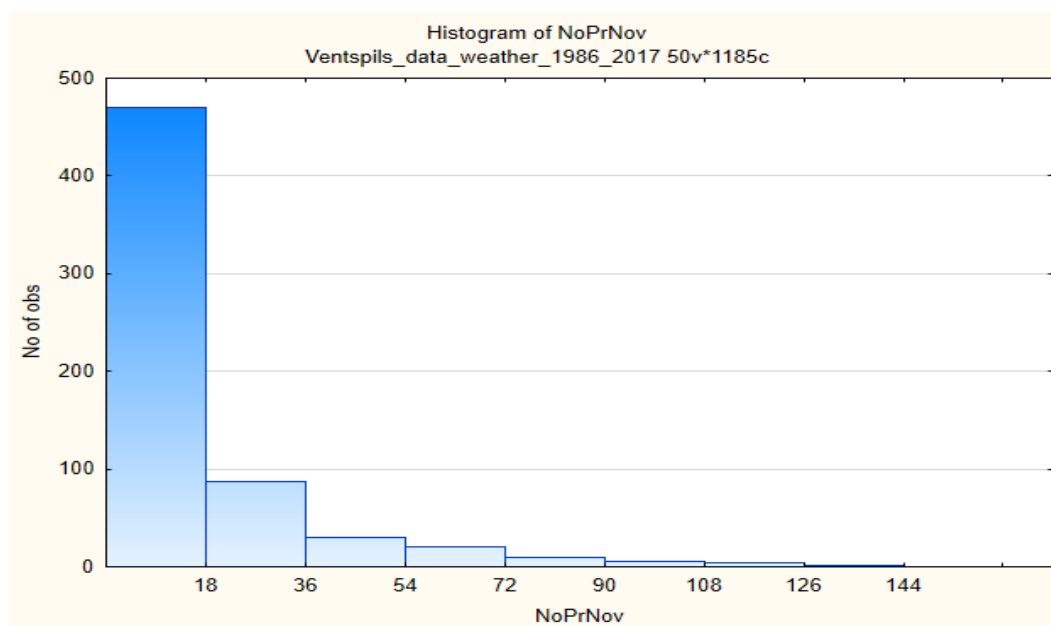


Fig.2. Histograms of the sojourn time in each state for each autumn month

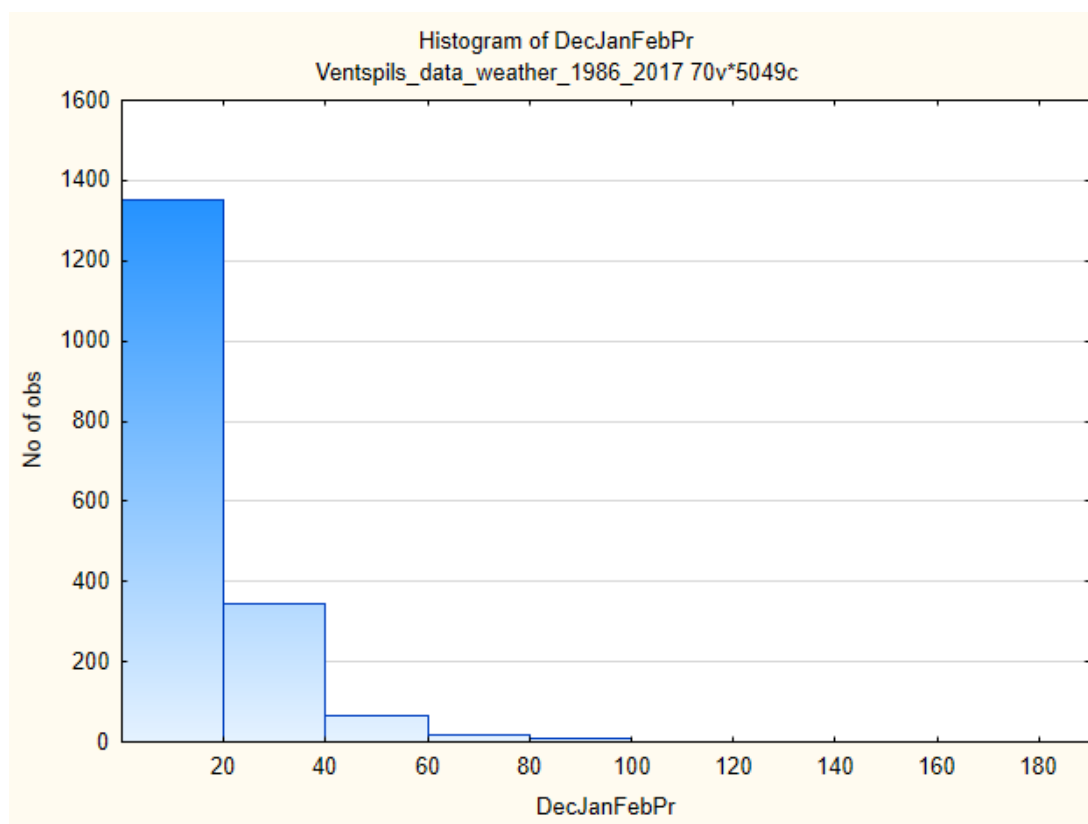
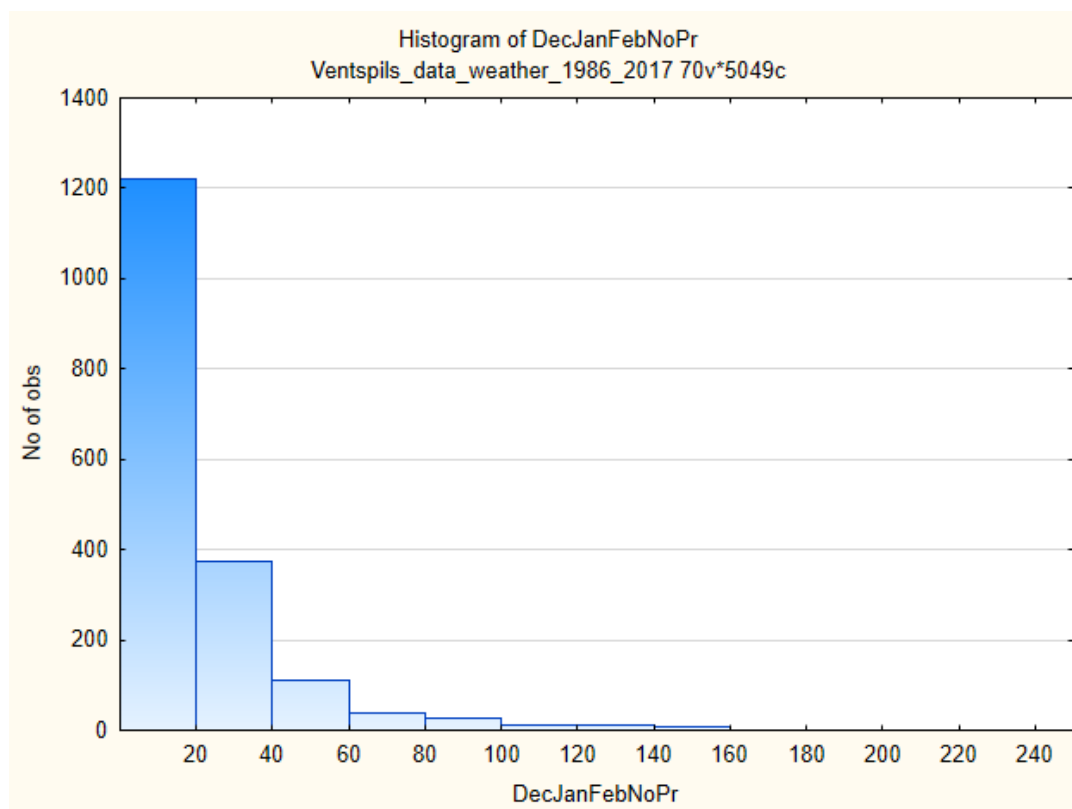


Fig.3. Histograms of the sojourn time in each state for all winter months combined